

# Adaptive Parameter Selection in Evolutionary Algorithms by Reinforcement Learning with Dynamic Discretization of Parameter Range

Arkady Rost  
ITMO University  
49 Kronverkskiy ave.  
Saint-Petersburg, Russia  
arkrost@gmail.com

Irina Petrova  
ITMO University  
49 Kronverkskiy ave.  
Saint-Petersburg, Russia  
petrova@rain.ifmo.ru

Arina Buzdalova  
ITMO University  
49 Kronverkskiy ave.  
Saint-Petersburg, Russia  
abuzdalova@gmail.com

## ABSTRACT

Online parameter controllers for evolutionary algorithms adjust values of parameters during the run of an evolutionary algorithm. Recently a new efficient parameter controller based on reinforcement learning was proposed by Karafotias et al. In this method ranges of parameters are discretized into several intervals before the run. However, performing adaptive discretization during the run may increase efficiency of an evolutionary algorithm. Aleti et al. proposed another efficient controller with adaptive discretization.

In the present paper we propose a parameter controller based on reinforcement learning with adaptive discretization. The proposed controller is compared with the existing parameter adjusting methods on several test problems using different configurations of an evolutionary algorithm. For the test problems, we consider four continuous functions, namely the sphere function, the Rosenbrock function, the Levi function and the Rastrigin function. Results show that the new controller outperforms the other controllers on most of the considered test problems.

## CCS Concepts

•Computing methodologies → Control methods; Reinforcement learning; Bio-inspired approaches; •Theory of computation → *Bio-inspired optimization*;

## Keywords

evolutionary algorithms; parameter control; reinforcement learning

## 1. INTRODUCTION

Let us denote efficiency of an evolutionary algorithm (EA) as the number of fitness function evaluations needed to find the optimal solution. Efficiency of EA is strongly correlated with its parameters. Common examples of such parameters are mutation and crossover probabilities. Optimal values of the parameters not only depend on the type of EA but also on the characteristics of the problem to be solved. Values of the parameters can be set before a run. However, optimal parameter values can change during a run, so an approach for adaptive parameter adjustment is required.

We consider parameters with continuous values. When adjusting such parameters, parameter ranges are usually discretized into some intervals. We can discretize parameter ranges a priori and keep the resulting segmentation during

a run. However, it was shown that adaptive discretization during optimization process may improve the performance of the algorithm [1, 2]. A possible explanation of this fact is as follows. The dynamic discretization allows to split the intervals into smaller subintervals. If the size of an interval is small, it is more likely to choose a good parameter value. Aleti et al. proposed entropy-based adaptive range parameter controller (EARPC) [1]. This method uses adaptive discretization.

Recently Karafotias et al. proposed another efficient parameter controller based on reinforcement learning [7, 8]. However, this method was not compared to EARPC. In the method proposed by Karafotias et al. a priori discretization is used. A new method which combines usage of reinforcement learning and dynamic discretization is proposed in this work.

The rest of the paper is organized as follows. First, the basic ideas used in EARPC and the method proposed by Karafotias are described. It is necessary to describe these ideas in order to explain how the proposed controller works. Second, two different versions of a new controller are proposed. Then the experiments are described. Finally, the new parameter controllers are compared with the other considered controllers.

## 2. RELATED WORK

Let us give the formal description of the adaptive parameter control problem. We have a set of  $n$  parameters  $v_1, \dots, v_n$ . The goal of the parameter controller is to select parameter values which maximize efficiency of EA. The first parameter controller to be considered is EARPC.

### 2.1 Entropy-based adaptive range parameter controller

The EARPC method [1] adjusts parameters separately. Let us denote a set of the selected parameter values as  $v = (v_1, \dots, v_n)$ . The efficiency of EA with parameters set to  $v$  is denoted as  $q(v)$ . During a run of the algorithm we save pairs of  $v$  and  $q(v)$ .

To select a new set of parameter values, we split the saved values of  $v$  into two clusters. For example, it can be done by k-means algorithm. Next, the range of each parameter  $v_i$  is split into two intervals. To select a split point, all saved values of  $v_i$  are sorted in ascending order. Candidates to be a split point are mid-points between two consequent values of  $v_i$  in the sorted list. For each candidate  $s_k$  to be

a split point, the set of the saved values of  $v_i$  is split into two subsets  $p_1$  and  $p_2$ . The subset  $p_1$  contains values which are less than or equal to  $s_k$ . The subset  $p_2$  contains values greater than  $s_k$ . Denote the number of the saved values of  $v_i$  which are contained in the cluster  $c_i$  and the subset  $p_j$  as  $c_i(p_j)$ . The entropy  $H$  is calculated according to Eq. 1 for each candidate  $s_k$ .

$$e_{p_1} = -\frac{|c_1(p_1)|}{|p_1|} \ln \left( \frac{|c_1(p_1)|}{|p_1|} \right) - \frac{|c_2(p_1)|}{|p_1|} \ln \left( \frac{|c_2(p_1)|}{|p_1|} \right),$$

$$e_{p_2} = -\frac{|c_1(p_2)|}{|p_2|} \ln \left( \frac{|c_1(p_2)|}{|p_2|} \right) - \frac{|c_2(p_2)|}{|p_2|} \ln \left( \frac{|c_2(p_2)|}{|p_2|} \right),$$

$$H = \frac{|p_1|}{|c_1|} e_{p_1} + \frac{|p_2|}{|c_2|} e_{p_2} \quad (1)$$

The split point  $s$  with the minimal value of entropy is selected. Two intervals  $[min, s]$  and  $(s, max]$  are obtained, where  $min$  and  $max$  are the minimal and the maximal possible values of  $v_i$  correspondingly. Values from the set  $p_1$  correspond to the first interval, values from the set  $p_2$  correspond to the second interval.

To decide from which interval we should choose a new value for the parameter  $v_i$ , we calculate the average quality of the parameter values in each interval. Let  $Q_1$  and  $Q_2$  denote the average quality of the values in the first and the second intervals correspondingly.  $Q_j$  is calculated according to Eq. 2.

$$Q_j = \frac{1}{|p_j|} \sum_{v \in p_j} q(v) \quad (2)$$

Then we randomly select interval, the probability of selection of interval  $j$  is proportional to  $Q_j$ . A new value of  $v_i$  is randomly chosen from the selected interval. The pseudocode of EARPC is presented in Algorithm 1. We wrote this pseudocode based on the description of EARPC given in [1].

---

**Algorithm 1** EARPC algorithm proposed by Aleti et al.

---

- 1: Earlier selected and saved values of  $v$  are split into two clusters  $c_1$  and  $c_2$
  - 2: **for** each parameter  $v_i$  **do**
  - 3:   Sort saved values of  $v_i$  in ascending order
  - 4:    $H_{best} \leftarrow \infty$
  - 5:   **for** split point  $s_k = \frac{v_{ij} + v_{i(j+1)}}{2}$  **do**
  - 6:     Split saved values of  $v_i$  into two sets  $p_1$  and  $p_2$  according to  $s_k$
  - 7:     Calculate entropy  $H$  according to Eq. 1
  - 8:     **if**  $H_{best} < H$  **then**
  - 9:        $H_{best} \leftarrow H$
  - 10:       $s \leftarrow s_k$
  - 11:   Split saved values of  $v_i$  into two sets  $p_1$  and  $p_2$  according to  $s$
  - 12:    $Q_1 = \frac{1}{|p_1|} \sum_{v \in p_1} q(v)$ ,  $Q_2 = \frac{1}{|p_2|} \sum_{v \in p_2} q(v)$
  - 13:   Randomly select interval, the probability of selection of interval  $[min, s]$  is proportional to  $Q_1$ , the probability of selection of interval  $(s, max]$  is proportional to  $Q_2$
  - 14:   Randomly select value of  $v_i$  from the selected interval
- 

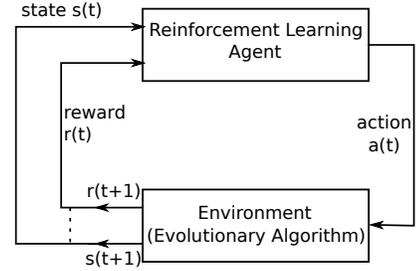


Figure 1: Reinforcement learning scheme

## 2.2 Parameter selection by reinforcement learning

Let us describe the general scheme of using reinforcement learning for parameter control in EA [3, 5, 10, 11]. In reinforcement learning (RL) [6, 12], the agent applies an action to the environment, then the environment returns some representation of its state and a numerical reward to the agent, and the process repeats. The goal of RL is to maximize the total reward [12].

While adjusting parameter of EA by RL, EA is treated as an environment. An action is selection of parameter values. EA generates new population using the selected parameter values. The obtained reward is based on the difference of the maximal fitness in two sequential iterations.

Let us describe how the agent selects parameter values. The parameter range is discretized into several intervals. Each interval corresponds to agent's action. To apply an action, the agent selects an interval and sets a random value from this interval as the parameter value. The method is illustrated in Fig. 1, where  $t$  is the number of the current iteration.

In the method proposed by Karafotias et al. [7, 8], a modification of the  $\varepsilon$ -greedy Q-learning algorithm is used. Let  $k$  denote the number of parameters being adjusted. The range of parameter  $v_i$  is discretized a priori into  $m_i$  intervals. An action of the agent consists of selection of intervals for each parameter and random selection of parameters values from the selected intervals. Thus the number of possible actions of the agent is  $\prod_{i=1}^k m_i$ . The obtained reward is calculated according to Eq. 3, where  $f_t$  is the best fitness function value obtained on the iteration  $t$  and  $c$  is a constant.

$$r = c \cdot \left( \frac{f_{t+1}}{f_t} - 1 \right) \quad (3)$$

Note that reward is always positive unless EA worsens the best obtained solution. To reduce the learning rate in the case of zero reward, the learning rate  $\alpha$  is changed to  $\alpha_0$  according to Eq. 4. Note that  $\alpha_0 \ll \alpha$ .

$$\alpha(r) = \begin{cases} \alpha, & \text{if } r > 0 \\ \alpha_0 \ll \alpha, & \text{in other cases} \end{cases} \quad (4)$$

To define the state of the environment, the following observables derived from the state of the EA are used [9]:

- genotypic diversity;
- phenotypic diversity (when different from genotypic);
- fitness standard deviation;

- stagnation counter (the number of iterations without fitness improvement);
- fitness improvement.

The dynamic state space segmentation method is used [13]. In this method, states are represented as a binary decision tree. Each internal node contains a condition on an observable. Leaf nodes represent environment states. For each state  $s$  array of  $Q(s, a)$  is stored, where  $Q(s, a)$  is the expected reward for each action in the state  $s$ . The expected reward in this state is denoted as  $V(s) = \max_a Q(s, a)$ .

Initially, the tree consists of one leaf which corresponds to a single state  $s$  and  $V(s) = 0$ . Each iteration of the algorithm consists of two phases: the data gathering phase and the processing phase. In the data gathering phase we obtain values of observables  $I = \{o_1 \dots o_m\}$  from EA, where  $m$  is the number of observables. Then we go down the tree using  $I$  and get the state  $s$  of the environment. The agent selects an action  $a$  using  $\varepsilon$ -greedy strategy, obtains reward  $r$  and new values  $I'$  of observables. Then the agent refreshes  $Q(s, a)$  and stores the resulting tuple  $(I, a, I', r)$ .

In the processing phase, we try to split the state  $s$  into two new states, which means that the leaf corresponding to the state  $s$  becomes an internal node and two new leaves are added as its children. To convert the leaf into a decision node, we have to choose an observable and a splitting value. For each saved tuple  $(I, a, I', r)$  in the leaf we calculate an estimated reward obtained after applying the action  $a$  to EA with the values  $I$  of observables. We denote this reward as  $q(I, a)$ . The value of  $q(I, a)$  is calculated as  $q(I, a) = r + \gamma V(s')$ , where  $s'$  corresponds to the values  $I'$  of observables, and  $\gamma$  is a constant called *discount factor*. For each observable  $o$  the tuples saved in the leaf are sorted in ascending order according to the value of  $o$  taken from  $I$ . Candidates to be chosen as a split point are mid-points between two consequent values of  $o$  in the sorted list. The saved values  $I$  taken from tuples stored in the leaf are divided into two subsets according to a candidate split point. We form two samples by dividing the calculated  $q(I, a)$  according to obtained subsets of  $I$ . A Kolmogorov-Smirnov criterion is run on these two samples and the obtained p-value is saved. After all split point candidates for all observables are checked, the smallest obtained p-value is selected. If it is smaller than 0.05, then the node is split at the corresponding observable and the corresponding split point.

The authors of the article where dynamic state space segmentation method was proposed [13] suggest that  $Q$  values should be recalculated after each split of a state. However, it is not obvious how to do it. So in the method proposed by Karafotias et al.  $Q$  and  $V$  values of two new nodes are set to the values of the parent node. The tuples saved in the parent node are split according to the chosen split point and the resulting parts are given to the corresponding children nodes. The pseudocode of the algorithm proposed by Karafotias et al. is presented in Algorithm 2. We wrote this pseudocode based on the description of this algorithm given in [7].

### 3. PROPOSED METHODS

We propose two new controllers. The first of them combines the EARPC algorithm and the approach proposed by Karafotias et al. The second one is based on reinforcement

---

**Algorithm 2** Algorithm proposed by Karafotias et al.

---

#### Data gathering phase

- 1: Initialize binary search tree: state  $s$ ,  $V(s) = 0$ ,  $Q(s, a) \leftarrow 0$  for each action  $a$
- 2: Obtain  $I$  — values of observables from EA
- 3: Go down the tree using  $I$  and find environment state  $s$
- 4: Select action  $a$  using  $\varepsilon$ -greedy strategy
- 5: Apply action  $a$  to environment, obtain reward  $r$
- 6: Obtain  $I'$  — values of observables from EA
- 7: Go down the tree using  $I'$  and find environment state  $s'$
- 8: Store tuple  $(I, a, I', r)$  in state  $s$
- 9:  $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$
- 10:  $V(s) = \max_a Q(s, a)$

#### Processing phase

- 1: **for** each tuple  $(I, a, I', r)$  in state  $s$  **do**
  - 2:   Calculate  $q(I, a) = r + \gamma V(s')$
  - 3:  $best \leftarrow +\infty$
  - 4: **for** observable  $o$  **do**
  - 5:   Sort tuples stored in state  $s$  according to value of  $o$  from  $I$ .
  - 6:    $tuples\_count \leftarrow$  number of tuples
  - 7:   **for**  $j \leftarrow 1$  to  $tuples\_count$  **do**
  - 8:      $I_j \leftarrow$   $j$ -th tuple in the sorted list
  - 9:      $o_j \leftarrow$  value of observable  $o$  in  $I_j$
  - 10:      $candidate \leftarrow \frac{o_j + o_{j+1}}{2}$
  - 11:      $x \leftarrow \{q(I_j, a) | o_j \leq candidate\}$
  - 12:      $y \leftarrow \{q(I_j, a) | o_j > candidate\}$
  - 13:     Calculate  $p$ -value for  $x$  and  $y$  using Kolmogorov-Smirnov criterion
  - 14:     **if**  $p$ -value  $<$   $best$  **then**
  - 15:        $best \leftarrow p$ -value
  - 16:        $best\_observable \leftarrow o$
  - 17:        $best\_split \leftarrow candidate$
  - 18: **if**  $best < 0.05$  **then**
  - 19:   Create two new states  $s_1$  and  $s_2$
  - 20:   Split tuples stored in  $s$  to  $s_1$  and  $s_2$  corresponding to  $best\_observable$  and  $best\_split$
  - 21:   Copy  $Q(s, a)$  into  $Q(s_1, a)$  and  $Q(s_2, a)$
  - 22:   Replace node corresponding to state  $s$  with internal node with two new leaves, corresponding to states  $s_1$  and  $s_2$
- 

learning and splitting of parameter ranges using Kolmogorov-Smirnov criterion.

### 3.1 Method which combines Karafotias et al. and EARPC

In the method proposed by Karafotias et al., dynamic state space segmentation is used. However, ranges of parameters being adjusted are discretized a priori. We propose to improve the method proposed by Karafotias et al. by using EARPC method for dynamic discretization of ranges of parameters. When we change discretization of the parameter range, we change the set of agent actions. So we have to change the process of selection of the parameter values in the method proposed by Karafotias et al.

The values of parameters are selected as in the EARPC method. Therefore, on each iteration of the algorithm we obtain values  $I$  of observables of EA and go down the binary

decision tree of states (see Section 2.2) to find a state  $s$  corresponding to  $I$ . Then we select values  $v$  of parameters, get new values of observables  $I'$  and save the tuple  $(I, v, I', r)$  in the state  $s$ , where  $v$  are selected values of parameters and  $r$  is calculated according to Eq. 3. The values  $v$  of parameters are selected by EARPC method using the tuples saved earlier in the state  $s$ . To apply the EARPC algorithm, we have to calculate  $q(v)$  which is the efficiency measure of EA with the parameters set to  $v$ . We use  $r$  from the tuple  $(I, v, I', r)$  as the efficiency measure  $q(v)$ .

To split the states, we have to calculate  $q(I, a) = r + \gamma V(s')$ , where  $s'$  corresponds to  $I'$ . In the proposed method actions of the agent are changing during the run. So we cannot calculate  $V(s')$  as  $\max_a Q(s', a)$ . In this case we use the expected value of the reward obtained in the state  $s'$  as  $V(s')$ . The EARPC algorithm splits the parameter range into two intervals. One of these intervals is selected with probability proportional to average reward on this interval. So  $V(s')$  is calculated as  $\sum_{i=1}^2 \frac{Q_i^2}{Q_1 + Q_2}$ , where  $Q_1$  and  $Q_2$  are average rewards on the two intervals.

To the best of our knowledge, there is no specific method of recalculating  $Q$ -values after splitting a state. In the method proposed by Karafotias et al.,  $Q$  and  $V$  values of two new nodes are set to the values of the parent node. In this method we do not store  $Q$ -values in leafs. So we do not need to recalculate  $Q$  and  $V$  for two new states after splitting.

### 3.2 Method with adaptive selection of action set

Preliminary experiments showed that there was no significant improvement of the EA efficiency when high number of states was used. Thus the second proposed method described in this section does not use the binary decision tree of states, although it is used in the first proposed method and the method proposed by Karafotias et al. In the second proposed method we have a single state of the environment. In the method proposed by Karafotias et al., values of all parameters to be adjusted are set simultaneously by the  $Q$ -learning agent. In the second proposed method the value of each parameter is set independently of the other parameters. So we have a separate  $Q$ -learning agent for each parameter.

Initially, for each parameter  $v_i$  we have only one action of the agent, which corresponds to the selection of parameter value from the range  $[min, max]$ , where  $min$  and  $max$  are the minimal and the maximal possible values of the parameter  $v_i$ . Each agent $_i$  applies this action and sets the value of  $v_i$ . Then we use the selected values  $v = (v_1 \dots v_n)$  in EA and calculate the reward  $r$  according to Eq. 3. We store the tuple of the selected values  $v$  and the obtained reward  $r$ . These tuples are used for splitting the parameter range and, as a consequence, they are used for changing sets of actions of the agents. The agents apply the single possible action until enough tuples for splitting of the range are stored.

After enough tuples are stored, we search a split point for the range of each parameter using Kolmogorov-Smirnov criterion. The criterion is used in the same way as it was used by Karafotias et al. when splitting states. If the split point is not found, we do not split the range. If the split point is found, we obtain two intervals  $L$  and  $R$ . Then we try to split  $L$  and  $R$  in the same way. We repeat this process  $i$  times. So the maximum number of the intervals is  $2^i$ .

An agent selects an action using  $\varepsilon$ -greedy strategy until  $Q$ -values become almost equal for all actions. In this case the expected rewards for all actions are almost equal, therefore the agent can not select which action is the most efficient. So the range of the parameter is re-discretized. The pseudocode of the proposed method for the case when  $i = 2$  is presented in Algorithm 3.

---

**Algorithm 3** Algorithm with adaptive selection of action set

---

```

1: State  $s \leftarrow single\_state$ 
2: for each parameter  $v_i$  to be adjusted do
3:    $P_i \leftarrow \{[v_i^{min}, v_i^{max}]\}$ 
4:    $A_i \leftarrow$  actions corresponding to partition  $P_i$ , where  $A_i$ 
     is a set of actions for agent $_i$  adjusting parameter  $v_i$ 
5:    $Q_i(s, a) \leftarrow 0$ 
6:   for each parameter  $v_i$  to be adjusted do
7:     if  $P_i = \{[v_i^{min}, v_i^{max}]\}$  then
8:       Split of range ( $v_i$ )
9:     else if  $A_i$  contains two or more actions and  $Q(s, a) -$ 
      $Q(s, a') < \delta$  then
10:      Split of range ( $v_i$ )
11:     Agent $_i$  selects action  $a_i$  from  $A_i$  using  $\varepsilon$ -greedy strat-
     egy
12:   Apply actions  $a_1 \dots a_n$  to environment, obtain reward  $r$ 

13: for each selected action  $a_i$  do
14:    $Q(s, a_i) \leftarrow Q(s, a_i) + \alpha(r + \gamma \max_{a'_i} Q(s, a'_i) - Q(s, a_i))$ 

```

**Split of range of  $v_i$**

```

1: Sort saved tuples of  $(v, r)$  according to  $v_i$ 
2:  $tuples\_count \leftarrow$  number of tuples
3:  $best \leftarrow +\infty$ 
4: for  $j \leftarrow 1$  to  $tuples\_count$  do
5:    $v_{i,j} \leftarrow$  value of  $v_i$  in  $j$ -th tuple in sorted list
6:    $candidate \leftarrow \frac{v_{i,j} + v_{i,j+1}}{2}$ 
7:   for  $j \leftarrow 1$  to  $tuples\_count$  do
8:      $x \leftarrow \{r | v_{i,j} \leq candidate\}$ 
9:      $y \leftarrow \{r | v_{i,j} > candidate\}$ 
10:  Calculate  $p$ -value for  $x$  and  $y$  using Kolmogorov-
     Smirnov criterion
11:  if  $p$ -value  $< best$  then
12:     $best \leftarrow p$ -value
13:     $best\_split \leftarrow candidate$ 
14: if  $best > 0.05$  then
15:  Split saved tuples into sets  $L$  and  $R$  according to
      $best\_split$ 
16:  Find split point  $s_l$  for  $L$ 
17:  Find split point  $s_r$  for  $R$ 
18:  if Split points  $s_l$  and  $s_r$  are not found then
19:     $P_i \leftarrow \{[v_{min}, s], (s, v_{max})\}$ 
20:  else if split point  $s_l$  is not found then
21:     $P_i \leftarrow \{[v_{min}, s], (s, s_r], (s_r, v_{max})\}$ 
22:  else if split point  $s_r$  is not found then
23:     $P_i \leftarrow \{[v_{min}, s_l], (s_l, s], (s, v_{max})\}$ 
24:  else
25:     $P_i \leftarrow \{[v_{min}, s_l], (s_l, s], (s, s_r], (s_r, v_{max})\}$ 
26:  Set  $A_i \leftarrow$  actions corresponding to partition  $P_i$ 

```

---

## 4. EXPERIMENTS AND RESULTS

The proposed methods were compared with the EARPC

algorithm, the approach proposed by Karafotias et al. and the  $Q$ -learning algorithm. In the  $Q$ -learning algorithm a single state is used and the ranges of parameter values are discretized a priori on five equally sized intervals as in the algorithm proposed by Karafotias et al. The considered methods were tested on several real-valued functions with different landscapes and different number of local optima. We implemented the EARPC algorithm ourselves and we used the implementation of the method proposed by Karafotias et al. kindly given by the authors of [7].

## 4.1 Experiment description

Let us denote the optimized function as  $F(x_1, \dots, x_n) : R^n \rightarrow R$ ,  $x_i \in [min_i, max_i]$ . Then the goal of EA is to find a vector  $x_1 \dots x_n$ , such as the global minimum of the function with  $\epsilon$  precision is reached on this vector. The algorithms were tested on the sphere function (Eq. 5), the Rosenbrock function (Eq. 6), the Levi function (Eq. 7) and the Rastrigin function (Eq. 8).

$$f(x_1, \dots, x_n) = \sum_{i=1}^n x_i^2 \quad (5)$$

$$f(x_1, x_2) = (a - x_1^2)^2 + b(x_2 - x_1^2)^2 \quad (6)$$

$$f(x_1, x_2) = \sin^2(3\pi x) + (x - 1)^2(1 + \sin^2(3\pi y)) + (y - 1)^2(1 + \sin^2(2\pi y)) \quad (7)$$

$$f(x_1, \dots, x_n) = A \cdot n + \sum_{i=1}^n [x_i^2 - A \cos(2\pi x_i)] \quad (8)$$

An individual of EA is a real vector  $x_1 \dots x_n$ . Mutation operator is defined as in the work by Karafotias et al. For each  $x_i$  in a vector we apply the following transformation:

$$x_i = \begin{cases} max_i, & \text{if } x_i + \sigma dx_i > max_i \\ min_i, & \text{if } x_i + \sigma dx_i < min_i \\ x_i + \sigma dx_i, & \text{in other cases} \end{cases} \quad (9)$$

where  $1 \leq i \leq n$ ,  $dx_i \sim \mathcal{N}(0, 1)$ , and  $\sigma$  is the parameter to be adjusted. We expect that the closer to the global optimum EA is, the smaller  $\sigma$  should be. The range of  $\sigma$  is  $[0, k]$ , where  $k$  is a constant. As the range of  $\sigma$  grows, it becomes harder to find the optimal value of  $\sigma$ . We used different values of  $k$  presented in Table 1. Note that the higher  $k$  is, the greater range of  $\sigma$  is.

We used  $(\mu + \lambda)$  evolution strategy with different values of  $\mu$  and  $\lambda$  presented in Table 1. All the considered algorithms were run 30 times on each problem instance, then the results were averaged. Parameters of EA are presented in Table 1.

The reward function in reinforcement learning was defined as follows:  $r = c \cdot \frac{f_t - f_{t+1}}{f_{t+1}}$ , where  $f_t$  is the minimal fitness function value in the generation  $t$  and  $c$  is a constant taken from [7]. Value of  $c$  is presented in Table 2. Note that when we solve minimization problem with elitist selection,  $f_{t+1}$  is less than or equal to  $f_t$ , so the reward is always positive. Parameters of reinforcement learning were taken from [7]. They are presented in Table 2.

## 4.2 Results and discussion

The results are presented as follows. First, we present and analyze the obtained average number of generations needed

Table 1: EA parameters

Parameter	Description	Value
$k$	maximal value of $\sigma$	{1, 2, 3}
$\mu$	the number of parents	{1, 5, 10}
$\lambda$	the number of children	{1, 3, 7}
$\epsilon$	precision	$10^{-5}$

Table 2: RL parameters

Parameter	Description	Value
$\alpha$	learning rate	0.9
$\alpha_0$	learning rate in case of zero reward	0.02
$\gamma$	discount factor	0.8
$\epsilon$	exploration probability	0.1
$c$	coefficient in reward	100

to reach the optimum using the two proposed methods, the EARPC algorithm, the approach proposed by Karafotias et al. and the  $Q$ -learning algorithm. Second, we analyze the selected values of the adjusted parameter. Also for the algorithm with dynamic discretization of the parameter range we analyze values of selected split points.

### 4.2.1 Influence on efficiency of EA

The average number of generations needed to reach the optimum using different parameter controllers is presented in Table 3. The first three columns contain values of EA parameters  $k$ ,  $\mu$  and  $\lambda$  correspondingly. The next 20 columns contain results of optimizing the following functions: the sphere function, the Rastrigin function, the Levi function and the Rosenbrock function. In each run we adjust  $\sigma$ , the parameter of mutation from Eq. 9. We expect that the closer to the global optimum EA is, the smaller  $\sigma$  should be. For each function, we present the results of the following parameter controllers: the proposed method with adaptive selection of action set (A), the  $Q$ -learning algorithm (Q), the approach proposed by Karafotias et al. (K), the EARPC algorithm (E) and the proposed method which combines Karafotias et al. and EARPC (E+K). In the last row, for each algorithm the total number of the EA configurations on which this algorithm outperformed all other algorithms is presented.

For each problem instance, the average number of generations needed to reach the optimal value is presented. The gray background corresponds to the best result for each EA configuration. Note that we do not compare EA algorithms characterized by different values of  $\mu$  and  $\lambda$ . We only compare methods of parameter control on the same EA configuration. For each set of  $k$ ,  $\mu$ ,  $\lambda$  values the number of fitness function evaluations in a generation is the same for all parameter control methods. So using the number of generations instead of the number of fitness evaluations does not affect results of comparison within a row.

The deviation of the average number of generations needed to reach the optimum in the proposed method with adaptive selection of action set is less than 5%. For the  $Q$ -learning algorithm and the approach proposed by Karafotias et al. the deviation is about 10%. For the EARPC algorithm and the proposed method which combines Karafotias et al. and EARPC the deviation is about 40%.

Table 3: Averaged number of runs needed to reach the optimum using the the proposed method with adaptive selection of action set (A), the Q-learning algorithm (Q), the approach proposed by Karafotias et al. (K), the EARPC algorithm (E) and the proposed method which combines Karafotias et al. and EARPC (E+K)

		Sphere function					Rastrigin function					Levi function					Rosenbrock function					
$k$	$\mu$	$\lambda$	A	Q	K	E	E+K	A	Q	K	E	E+K	A	Q	K	E	E+K	A	Q	K	E	E+K
1	1	1	2434	8769	8048	5258	4830	3631	9653	8385	7794	8689	3496	7200	7265	7986	14092	5124	15058	13418	9003	12905
1	1	3	2207	4221	2683	3942	3070	1776	2226	2069	3148	3610	1980	3305	2903	2789	3688	2301	3914	4167	3553	2701
1	1	7	878	1085	2620	1653	2247	1226	1605	1757	1422	1422	1321	1584	1600	1923	3820	1411	1791	2330	2296	1941
1	5	1	1450	1664	2076	4472	3893	1666	1706	2281	4341	4530	1778	1898	1865	2372	2246	1859	2311	2809	5730	4453
1	5	3	569	706	824	1589	845	918	894	942	1958	2197	855	862	794	1632	2171	1053	869	748	1393	2749
1	5	7	368	390	473	642	568	502	570	675	1679	1329	356	606	444	1426	1236	401	533	665	1130	1258
1	10	1	703	959	747	1438	728	1008	1103	1105	2681	2926	884	840	804	2353	1451	1617	1497	1593	2213	3085
1	10	3	331	378	358	534	400	622	604	665	1460	1485	533	502	624	1491	1134	683	488	616	1161	1225
1	10	7	186	195	167	142	422	358	489	473	1103	1858	308	331	294	510	766	483	290	235	391	429
2	1	1	4342	25192	28523	31738	16182	4744	23663	21043	27865	19142	4947	13420	14789	25721	30653	5165	23371	27239	20005	20819
2	1	3	2333	7681	6478	5152	4233	1839	6405	6825	9748	8997	2020	6884	6601	7477	4612	3461	7295	7169	7166	13342
2	1	7	1464	3360	3739	1688	2956	1160	2388	2183	3806	4085	1205	3521	2370	2533	2967	1753	3488	2968	5874	7870
2	5	1	1891	3814	3468	4908	5173	2467	3944	4029	5961	5750	1935	3517	3369	7222	5365	1990	8076	5810	11153	11856
2	5	3	974	994	1163	1631	946	1128	1614	1780	2293	1720	1085	1231	1326	3039	2943	1396	2116	2705	3775	4542
2	5	7	616	716	756	1511	626	967	1012	1461	933	1180	770	792	1089	1153	2180	934	1001	1247	1775	2584
2	10	1	1164	2397	1833	778	876	1722	2108	2255	2411	2807	1803	1884	2197	5139	3072	2022	3415	2787	8768	4603
2	10	3	459	445	465	719	788	988	1420	1116	1486	1234	887	988	1006	1045	2381	1181	1037	1106	3719	1648
2	10	7	188	320	252	380	413	615	856	712	627	922	618	731	564	806	1064	707	811	666	1740	2558
3	1	1	8055	36698	29112	54868	30710	5159	29082	23305	24249	27327	5216	21470	21953	28900	35119	6535	28702	36597	32533	43832
3	1	3	2826	7845	9115	12446	11985	2821	7977	10726	6541	16040	2900	9029	7071	6966	10883	3391	11427	9873	13061	10068
3	1	7	1427	3907	4813	10118	9207	1517	4680	4266	4144	5408	1700	4139	4019	3375	2062	2573	5820	6462	9797	9775
3	5	1	2447	8328	5886	3348	12313	2704	6208	5419	5953	7665	3105	6966	6742	10480	9059	3148	12438	6791	8952	11581
3	5	3	1445	2790	2222	1379	1848	1544	2514	2096	2823	2304	1808	2467	2664	4593	3714	1721	2289	3673	4434	6799
3	5	7	896	919	777	1424	1105	902	1120	1629	1033	1055	1017	1929	1788	594	1439	1231	1626	1255	3951	4914
3	10	1	1996	2398	2531	3173	3745	2048	3747	3258	5095	3873	2071	2901	2943	5210	4814	2352	4473	6531	8320	10539
3	10	3	1053	1074	1206	1114	171	1296	1740	1523	1013	1028	1330	1462	1832	3140	2389	1677	2000	1650	2611	2479
3	10	7	409	391	427	410	537	955	995	1203	779	839	667	1117	799	521	845	1101	1137	1098	1742	1907
Summary			20	2	2	2	1	22	2	0	3	0	19	1	4	2	0	19	3	5	0	0

Overall, when  $k$  increases, the range of  $\sigma$  increases and selection of the optimal value of  $\sigma$  becomes harder. The proposed method with adaptive selection of action set outperformed all other considered methods on most problem instances. According to multiple sign test [4], the proposed method with adaptive selection of action set is distinguishable from the other methods at the level of statistical significance  $\alpha = 0.05$ .

Let us discuss the possible reasons why the proposed method with adaptive selection of action set outperformed the other considered methods. In the method proposed by Karafotias et al., in the case of zero reward, for some consequent iterations  $Q$ -values become zero. So the experience of the agent is lost. The proposed method with adaptive selection of action set does not have this drawback because experience is not saved only in  $Q$ -values.

In EARPC and the proposed method which combines Karafotias et al. and EARPC, the range of parameter  $\sigma$  is split into two almost equal intervals. The proposed method with adaptive selection of action set splits interval using another criterion which allows to split the range more precisely.

Consider the efficiency of using binary decision tree of states. The  $Q$ -learning algorithm and the algorithm proposed by Karafotias et al. demonstrate similar performance. The same is also true for the EARPC algorithm and the proposed method which combines Karafotias et al. and EARPC. Therefore, applying dynamic state space segmentation algorithms in parameter controllers does not seem to increase efficiency of EA when solving the considered problems.

#### 4.2.2 Selected parameter values and split points

In the Fig. 2 values of  $\sigma$  selected by the two proposed methods and the method proposed by Karafotias et al. during Rastrigin function optimization are presented. For the other considered functions plots are similar. For brevity, they are not presented. The horizontal axis refers to the number of an iteration, the vertical axis refers to the selected value of  $\sigma$ .

We can see that method which combines Karafotias et al. and EARPC (Fig. 2a) in the end of optimization selects  $\sigma$  almost randomly. The algorithm proposed by Karafotias et al. (Fig. 2c) continues selection of an action if it has achieved positive reward for application of this action. However, if the algorithm has not obtained positive reward in some consequent iterations, it loses information about efficient action selected earlier. We can see that the proposed method with adaptive selection of action set (Fig. 2b) in the beginning of the optimization selects  $\sigma$  randomly, but during the optimization process the selected  $\sigma$  converges to the optimal value.

For the two proposed methods, the selected split points are shown in Fig. 3. The algorithm proposed by Karafotias et al. discretizes the parameter range a priori, so the considered type of plot is not applicable. The split points selected by EARPC are not presented because EARPC and the first proposed method demonstrate similar performance. The horizontal axis in Fig. 3 refers to the number of an iteration, the vertical axis refers to the selected split points of the  $\sigma$  range.

We can see that in the proposed method which combines Karafotias et al. and EARPC (Fig. 3a), the parameter range

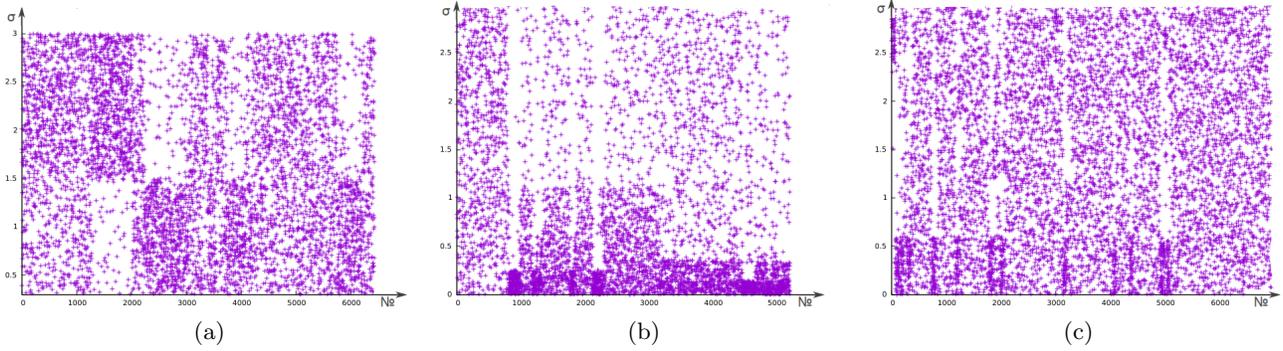


Figure 2: Selected values of  $\sigma$  in the proposed method which combines Karafotias et al. and EARPC (a), the proposed method with adaptive selection of action set (b) and the method proposed by Karafotias et al. (c) on Rastrigin function.

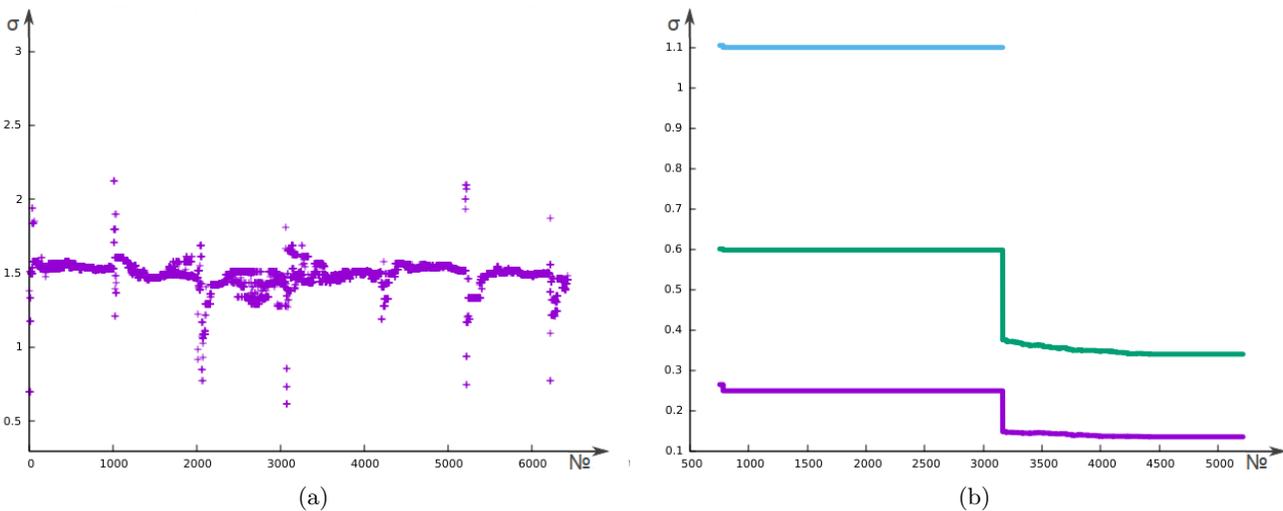


Figure 3: Split point values in the proposed method which combines Karafotias et al. and EARPC (a) and the proposed method with adaptive selection of action set (b) on Rastrigin function.

is split into two almost equal intervals on each iteration of EA. The proposed method with adaptive selection of action set (Fig. 3b) does not split interval until enough tuples of experience are obtained. Then (after about 750 iterations) the  $\sigma$  range is split into four intervals by three split points. Discretization of the  $\sigma$  range is not changed until  $Q$ -values become almost equal for all actions (after about 3150 iterations). Then the range is discretized again into three intervals. Note that the intervals after re-discretization are shrunk and the number of intervals is reduced. So the agent can select a good value of  $\sigma$  more precisely. This effect is reflected in Fig. 2b. After about 3150 iterations of the algorithm the agent almost always selects the value of  $\sigma$  close to the optimal value. The outliers can be explained by the fact that  $\varepsilon$ -greedy exploration strategy is used in  $Q$ -learning. According to this strategy, a random action is applied with probability  $\varepsilon$ .

## 5. CONCLUSION

We proposed two new parameter controllers based on reinforcement learning. These algorithms discretize parameter range dynamically. One of the proposed methods is based on two existing parameter controllers: EARPC and the algorithm proposed by Karafotias et al. In the second approach the parameter range is discretized using Kolmogorov-Smirnov criterion and it is re-discretized if the expected reward is almost equal for all actions of the agent.

The proposed methods were compared with EARPC, the algorithm proposed by Karafotias et al. and the  $Q$ -learning algorithm. We tested controllers with 27 configurations of EA on four test problems. On the most problem instances, the second proposed approach outperformed the other considered methods. We showed that this method improves the parameter value during the whole optimization process contrary to the other methods. It also may be noticed that application of dynamic state space segmentation algorithms in parameter controllers does not seem to increase efficiency of EA when solving the considered test problems.

## 6. REFERENCES

- [1] A. Aleti and I. Moser. Entropy-based adaptive range parameter control for evolutionary algorithms. In *Proceedings of Genetic and Evolutionary Computation Conference*, pages 1501–1508, 2013.
- [2] A. Aleti, I. Moser, and S. Mostaghim. Adaptive range parameter control. In *Proceedings of Congress on Evolutionary Computation*, pages 1–8, 2012.
- [3] F. Chen, Y. Gao, Z.-q. Chen, and S.-f. Chen. SCGA: Controlling genetic algorithms with sarsa (0). In *Computational Intelligence for Modelling, Control and Automation, 2005 and International Conference on Intelligent Agents, Web Technologies and Internet Commerce, International Conference on*, volume 1, pages 1177–1183. IEEE, 2005.
- [4] J. Derrac, S. Garcia, D. Molina, and F. Herrera. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm and Evolutionary Computation*, 1(1):3–18, 2011.
- [5] A. E. Eiben, M. Horvath, W. Kowalczyk, and M. C. Schut. Reinforcement Learning for Online Control of Evolutionary Algorithms. In *Proceedings of the 4th international conference on Engineering self-organising systems*, pages 151–160. Springer-Verlag, Berlin, Heidelberg, 2006.
- [6] A. Gosavi. Reinforcement Learning: A Tutorial Survey and Recent Advances. *INFORMS Journal on Computing*, 21(2):178–192, 2009.
- [7] G. Karafotias, Á. E. Eiben, and M. Hoogendoorn. Generic parameter control with reinforcement learning. In *Proceedings of Genetic and Evolutionary Computation Conference*, pages 1319–1326, 2014.
- [8] G. Karafotias, M. Hoogendoorn, and A. Eiben. Parameter control in evolutionary algorithms: Trends and challenges. *Evolutionary Computation, IEEE Transactions on*, PP(99):1–1, 2014.
- [9] G. Karafotias, S. K. Smit, and A. E. Eiben. A generic approach to parameter control. In *Proceedings of the 2012T European Conference on Applications of Evolutionary Computation, EvoApplications’12*, pages 366–375, Berlin, Heidelberg, 2012. Springer-Verlag.
- [10] J. E. Pettinger and R. M. Everson. Controlling Genetic Algorithms with Reinforcement Learning. In *Proceedings of Genetic and Evolutionary Computation Conference*, page 692, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc.
- [11] Y. Sakurai, K. Takada, T. Kawabe, and S. Tsuruta. A Method to Control Parameters of Evolutionary Algorithms by Using Reinforcement Learning. In *Proceedings of 2010 Sixth International Conference on Signal-Image Technology and Internet-Based Systems (SITIS)*, pages 74–79, 2010.
- [12] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, USA, 1998.
- [13] W. T. B. Uther and M. M. Veloso. Tree based discretization for continuous state space reinforcement learning. In *Proceedings of the Fifteenth National/Tenth Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence, AAAI ’98/IAAI ’98*, pages 769–774, Menlo Park, CA, USA, 1998.